

# 下一代RAG引擎 —— 技术挑战与实现

演讲人：张颖峰

InfiniFlow/ 创始人

# 极客邦科技 2024 年会议规划

促进软件开发及相关领域知识与创新的传播



# 目录

01 下一代RAG引擎

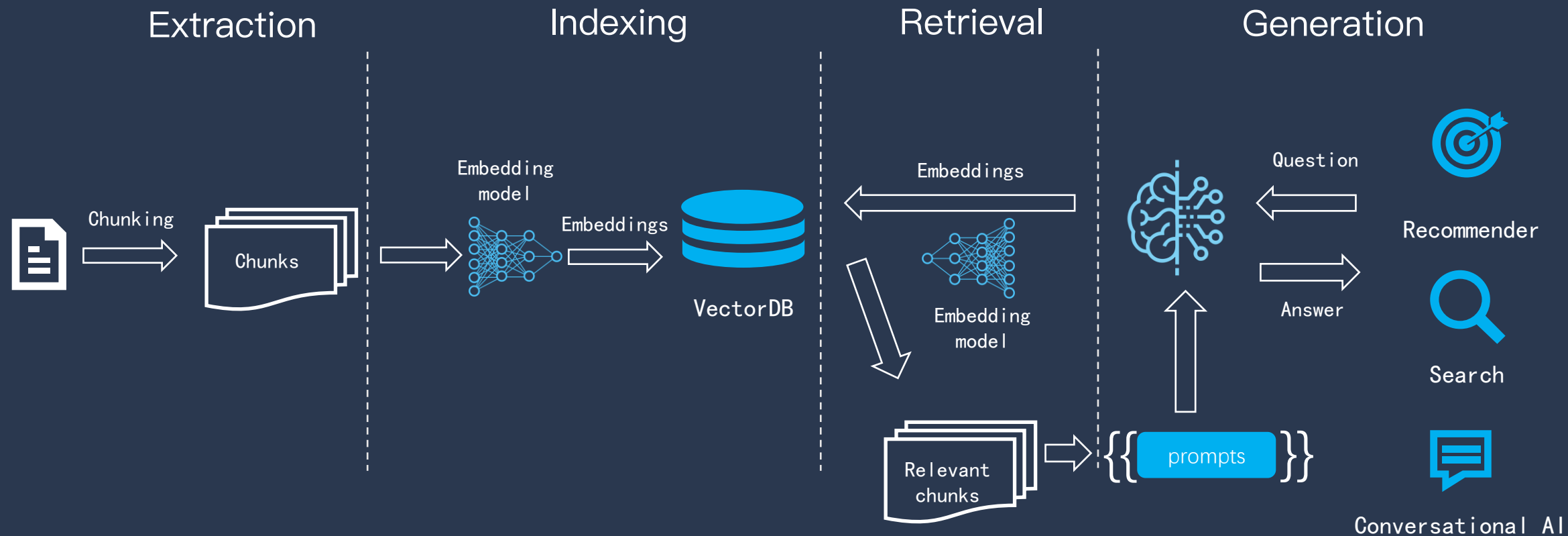
02 数据抽取模型

03 混合搜索

04 高级RAG

# 01 下一代RAG引擎

# RAG 架构模式

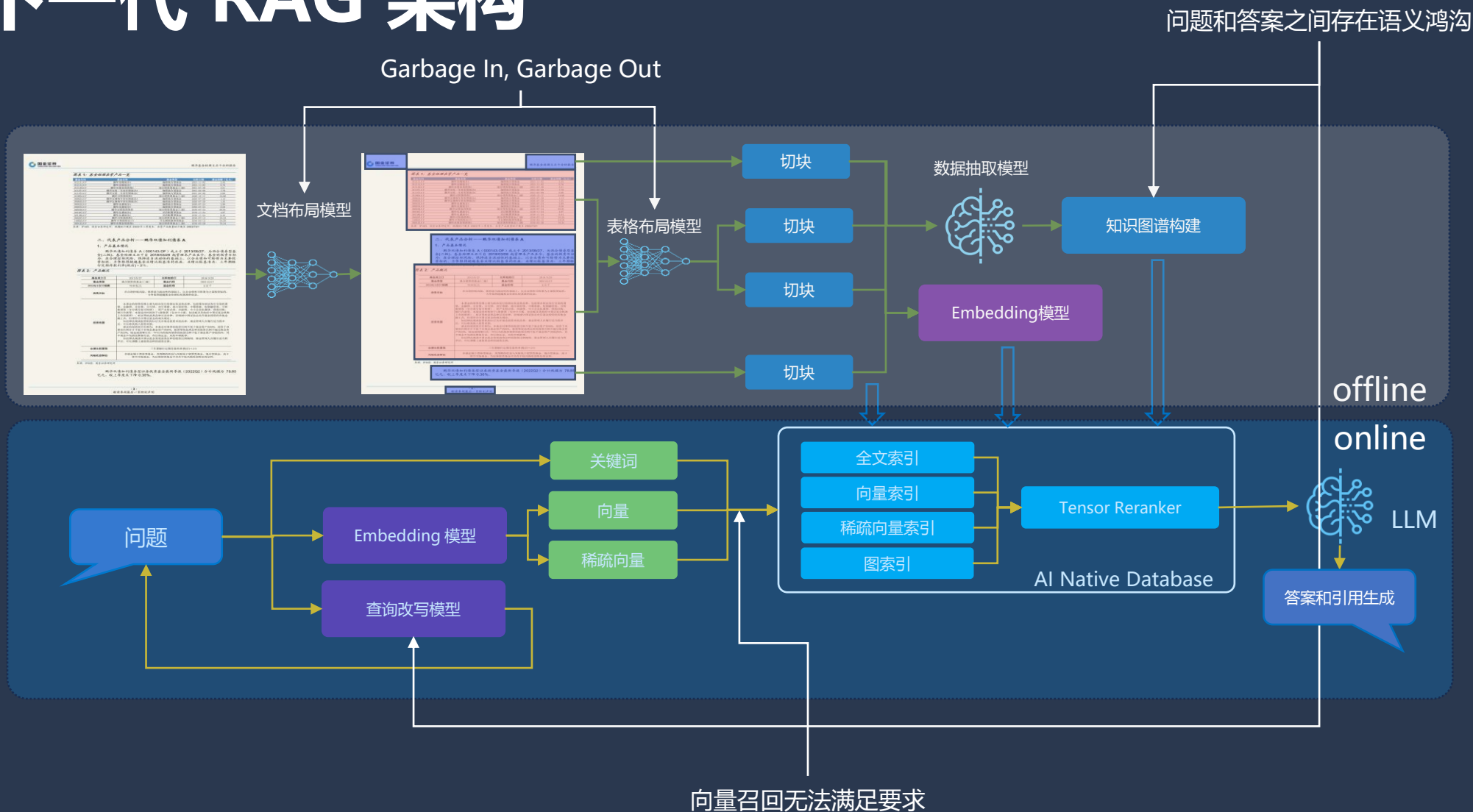


# ■ 当前 RAG 面临的挑战

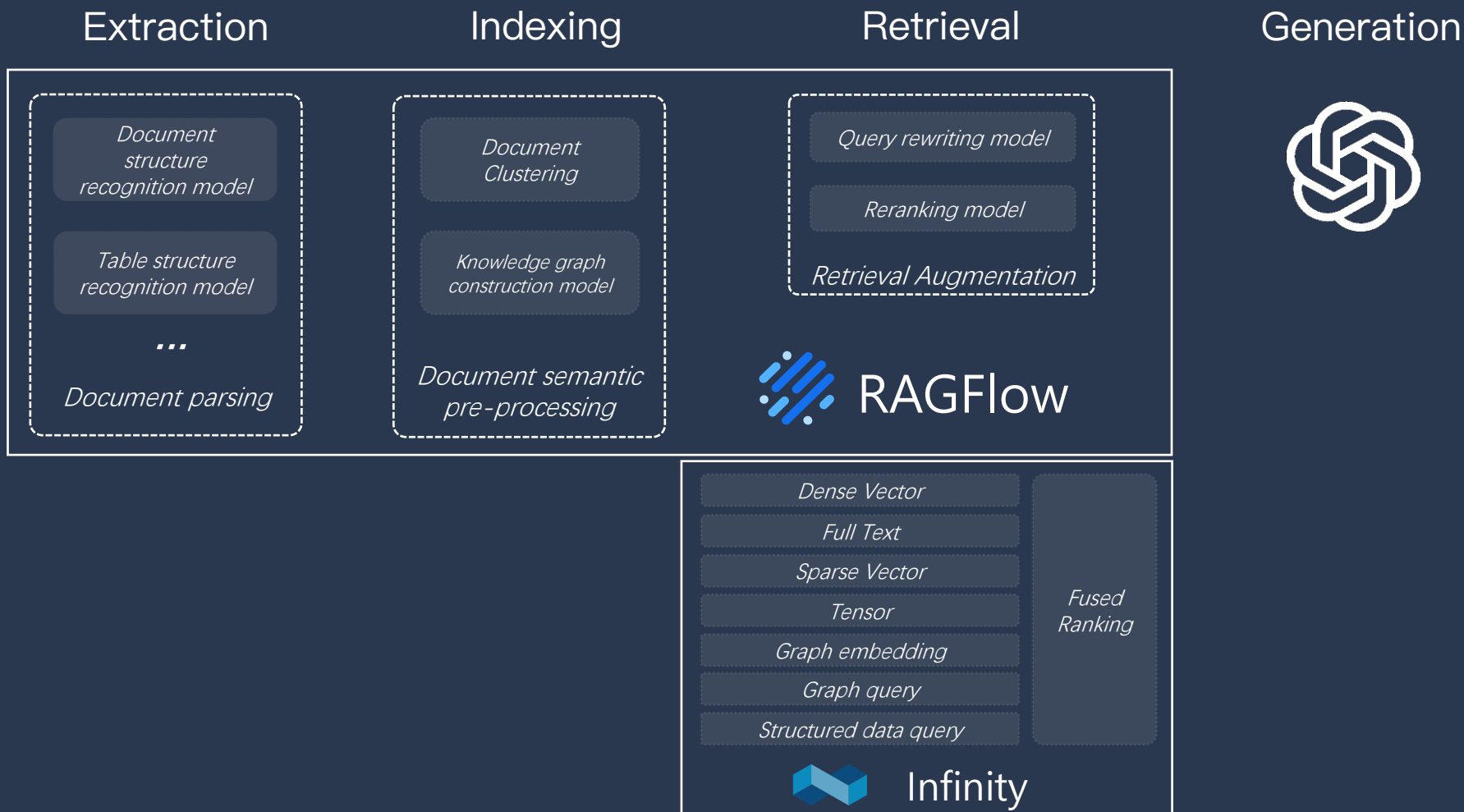
- 挑战一：向量的召回无法满足要求
- 挑战二：文档结构复杂，数据太乱， Garbage In, Garbage Out
- 挑战三：问题和答案所在文档关联不大，很难通过问题找到正确文档



# 下一代 RAG 架构



# Infinity + RAGFlow = Infiniflow





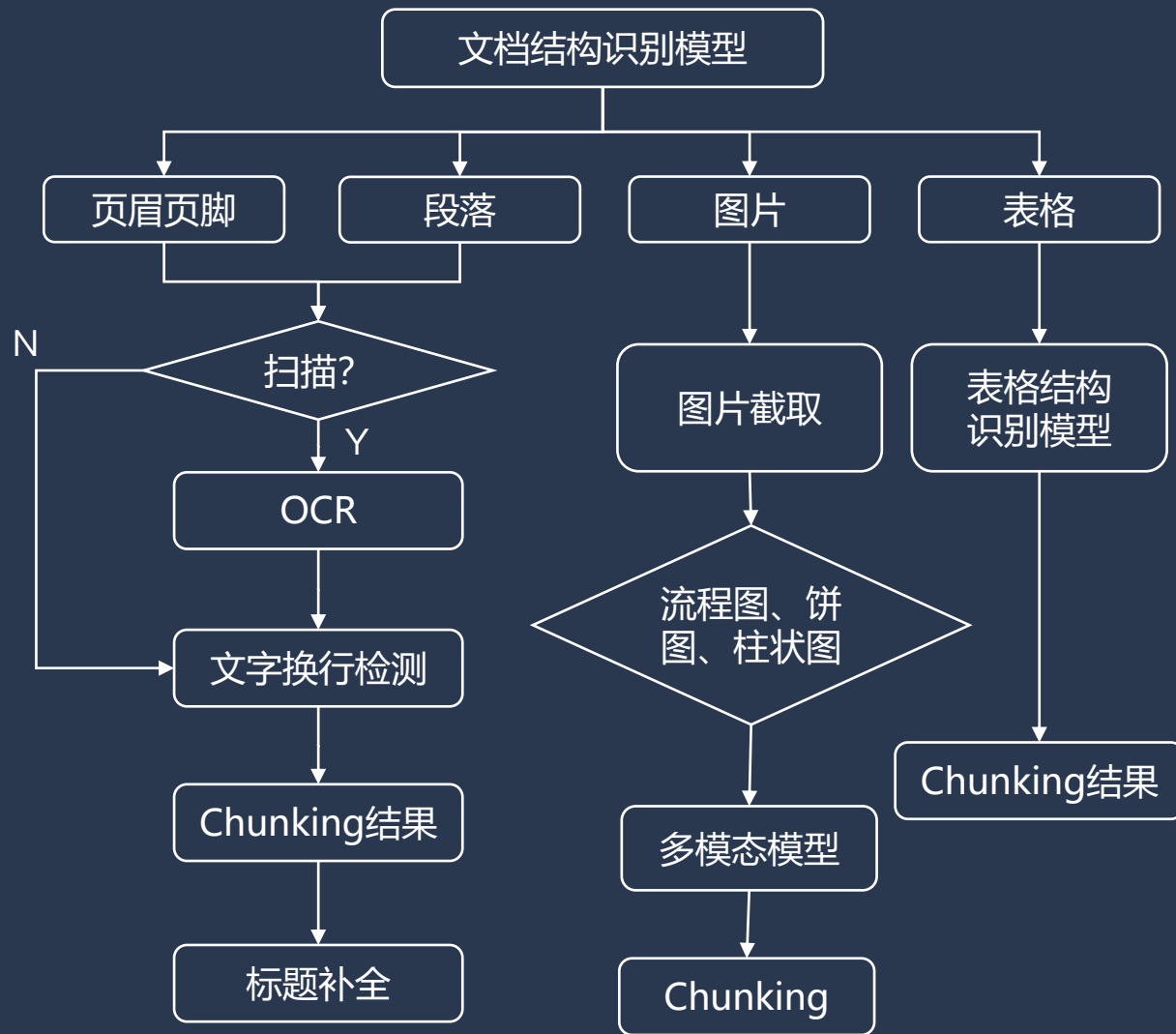
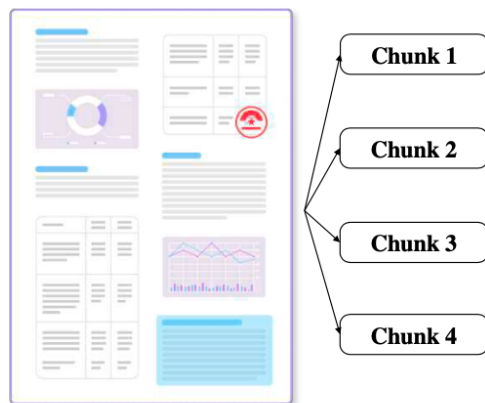
# 02 数据抽取模型

# 概要

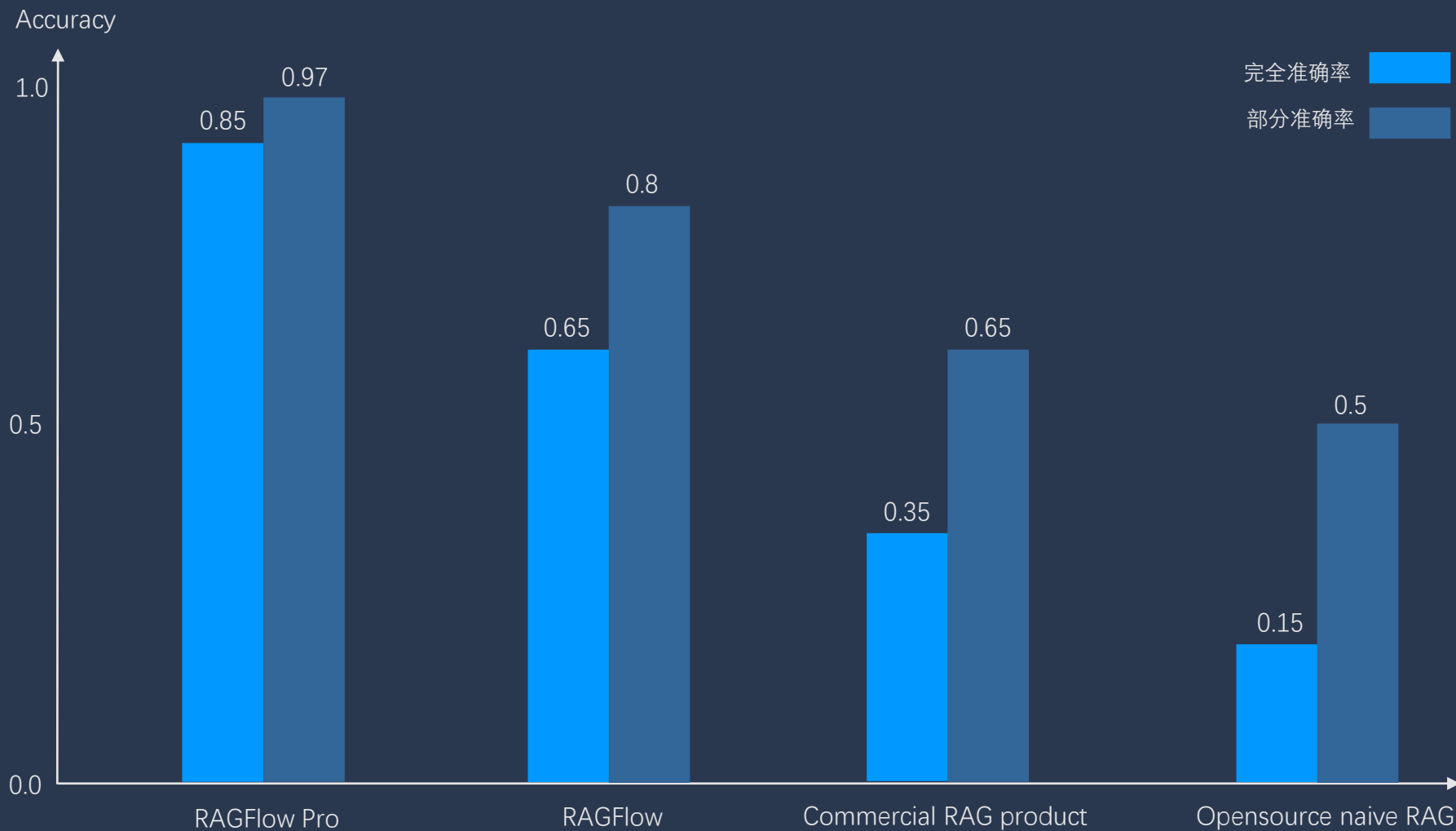


Documents

## Document Parsing & Chunking



# 调整抽取模型的 RAGFlow 对比



# 表格识别模型

板块	申万一级行业	2020Q 2	2020Q 4	2021Q 2	2021Q 4
制造	电力设备	14.37%	8.96%	15.63%	7.80%
	机械设备	3.03%	4.13%	4.54%	1.03%
	汽车	1.68%	2.45%	4.18%	1.33%
	公用事业	0.00%	0.00%	0.00%	5.32%
	环保	1.98%	0.96%	0.00%	2.67%
	国防军工	0.00%	7.42%	2.16%	6.08%
	综合	0.00%	0.00%	0.00%	0.00%
消费	纺织服饰	0.00%	0.00%	0.00%	0.00%
	家用电器	0.00%	1.06%	0.00%	0.90%
	农林牧渔	3.56%	0.00%	4.68%	4.10%
	轻工制造	2.95%	2.06%	0.86%	0.00%
	商贸零售	2.69%	1.25%	3.64%	0.15%
	食品饮料	7.90%	4.22%	0.83%	10.08%
	社会服务	1.53%	0.30%	1.03%	0.00%
医药	美容护理	0.00%	0.00%	2.20%	0.35%
	医药生物	15.64%	12.52%	9.55%	5.98%
科技	传媒	3.44%	0.00%	0.00%	2.50%
	电子	11.32%	9.25%	17.68%	11.48%
	计算机	11.21%	8.33%	2.62%	1.80%
	通信	2.70%	2.08%	0.00%	0.00%
金融	房地产	0.00%	0.00%	0.00%	0.70%
	非银金融	2.88%	8.31%	2.23%	9.34%
	银行	0.00%	4.18%	9.78%	0.36%
周期	煤炭	0.00%	0.00%	2.57%	0.00%
	石油石化	0.00%	1.15%	3.00%	0.00%
	钢铁	0.00%	3.45%	3.13%	3.46%
	基础化工	5.56%	3.37%	2.64%	9.78%
	建筑材料	3.09%	5.49%	4.20%	0.00%
	有色金属	2.50%	5.10%	2.09%	4.37%
	建筑装饰	0.00%	1.30%	0.00%	3.63%
	交通运输	1.97%	2.68%	0.84%	6.79%

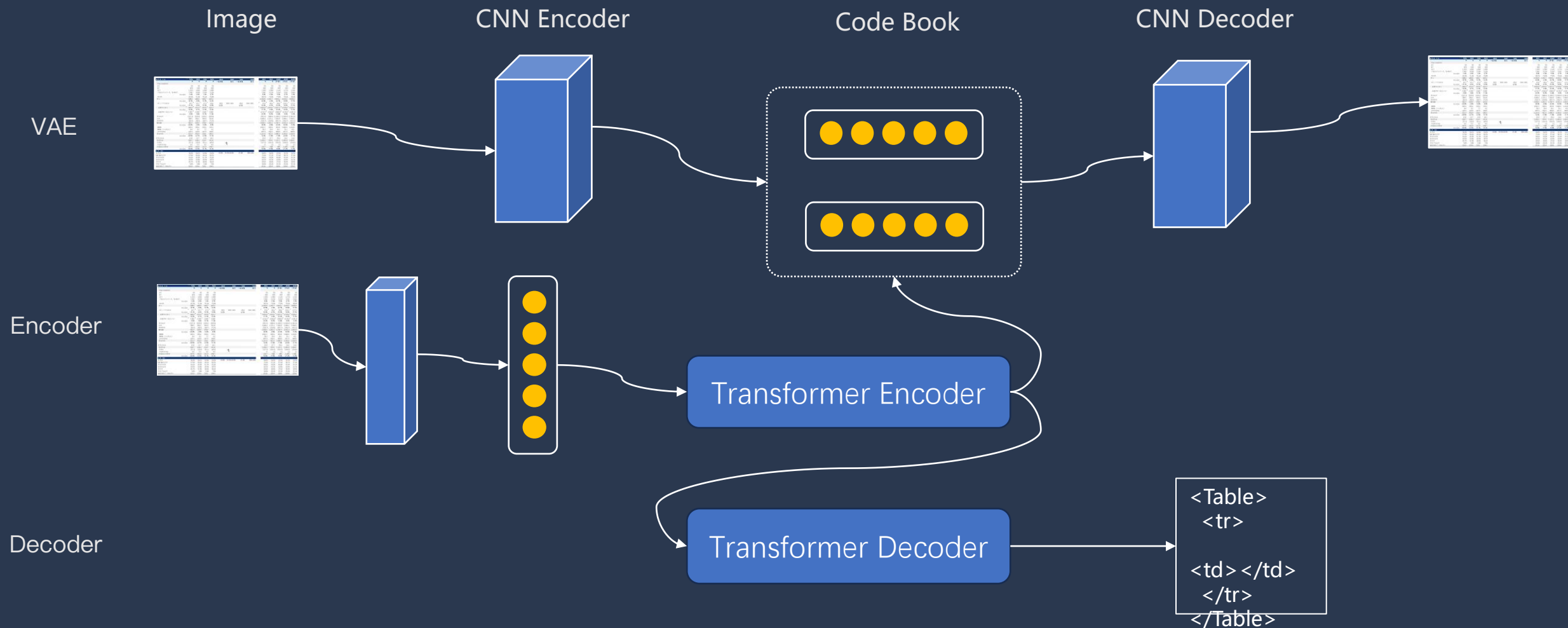
■ 单元格边界判定

■ 表头信息判定

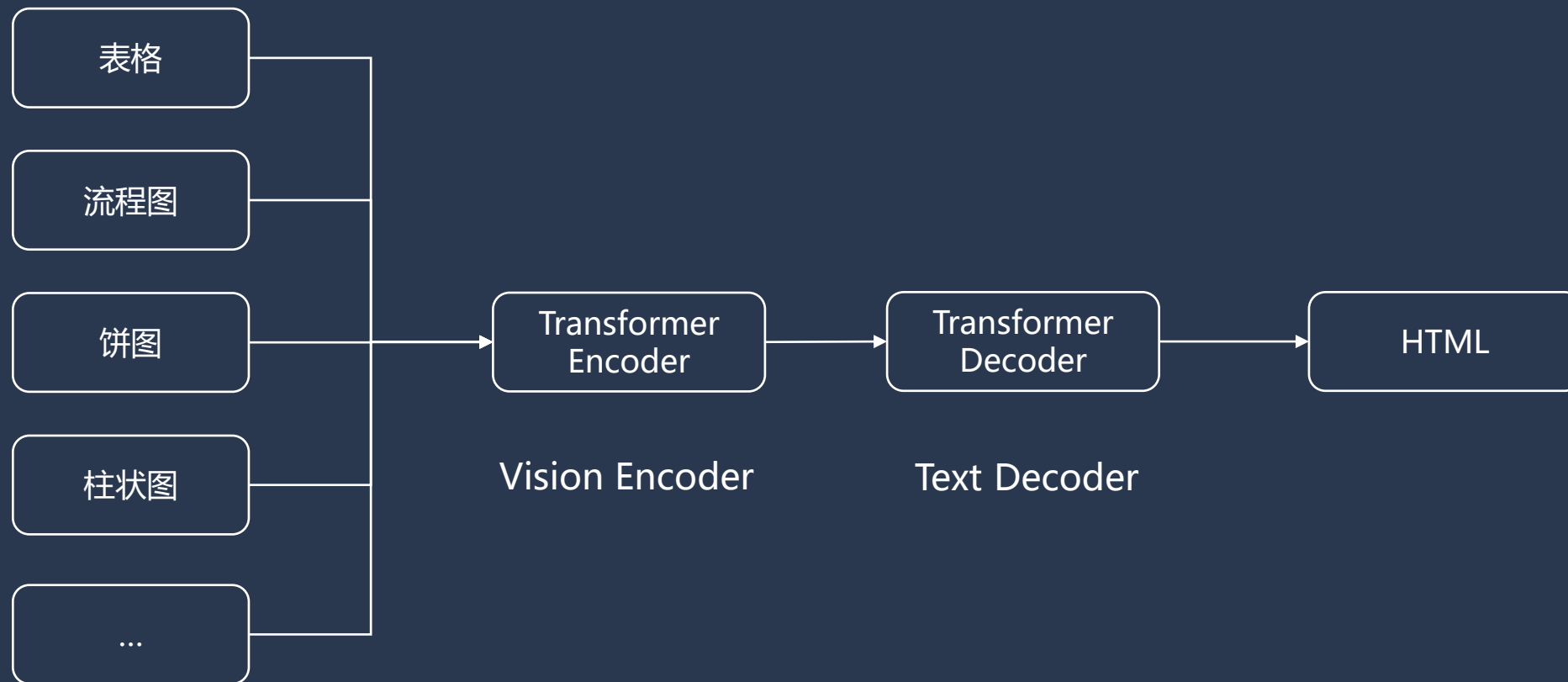
■ 单元格合并判定

■ 表格跨页判定

# 表格识别模型

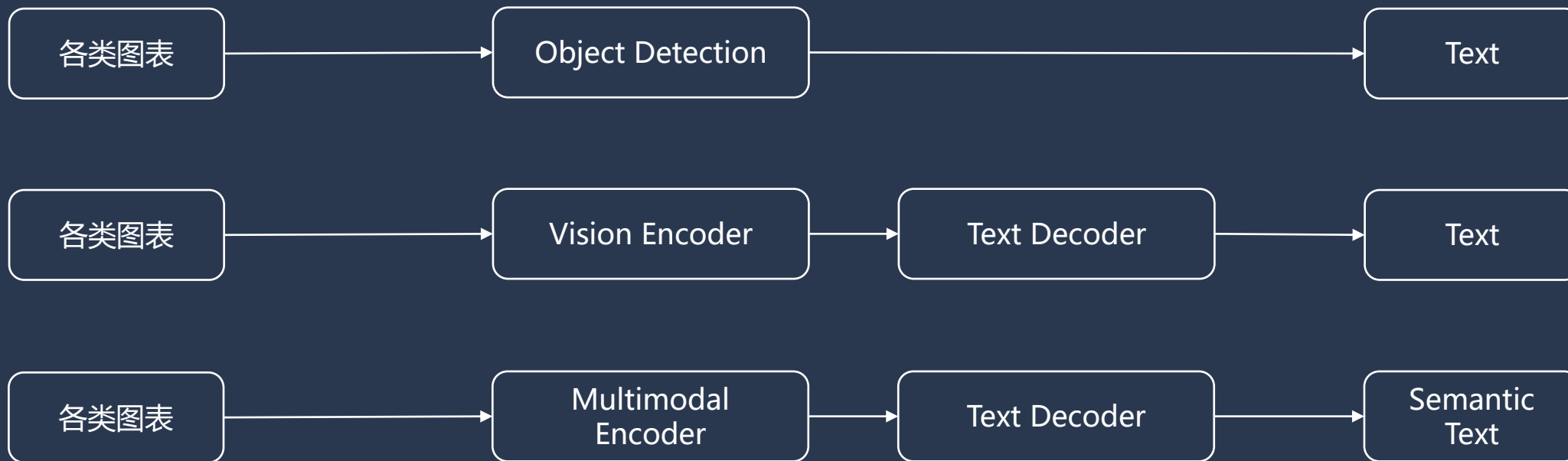


# 文档“大”模型





# “雕花” 还是多模态LLM?



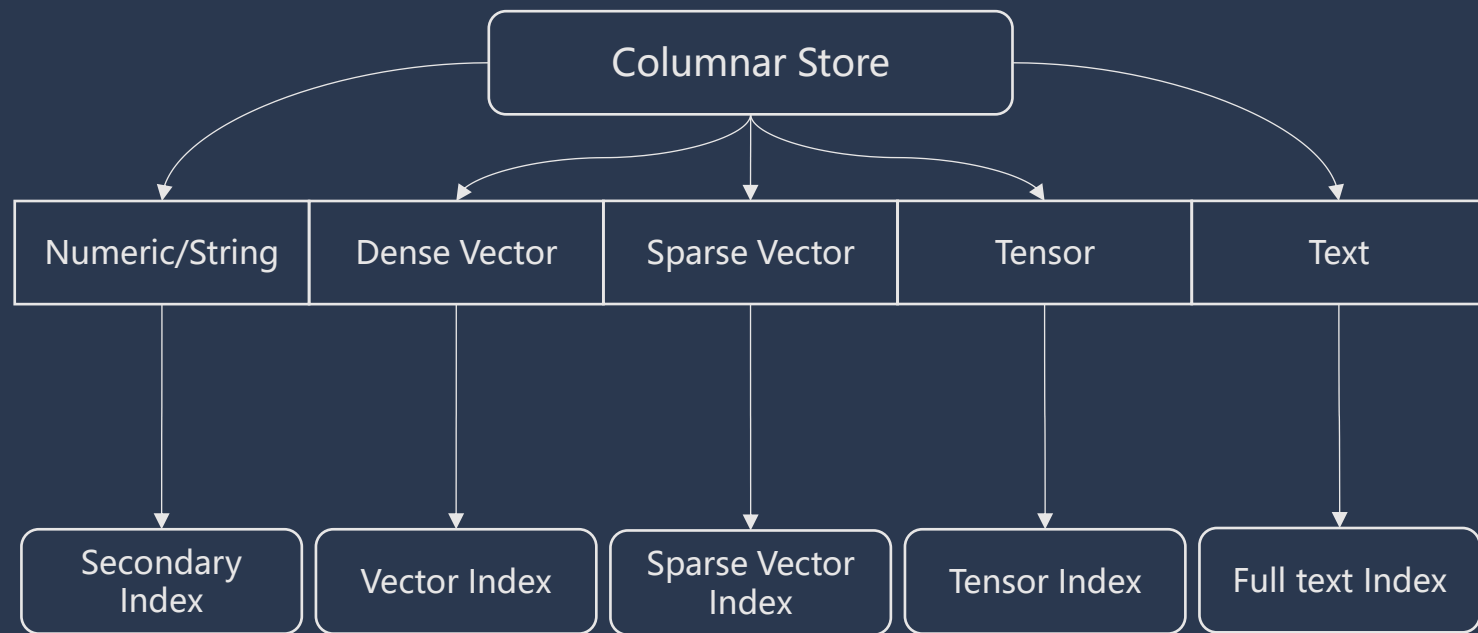
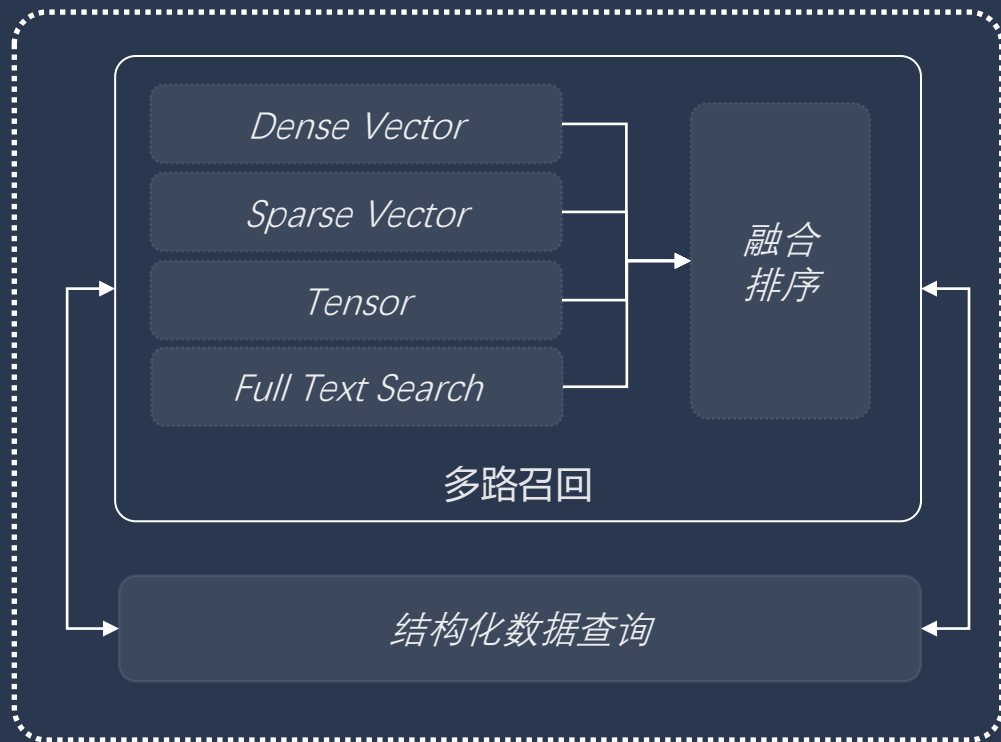
 **CharXiv: Charting Gaps in Realistic Chart Understanding in Multimodal LLMs**



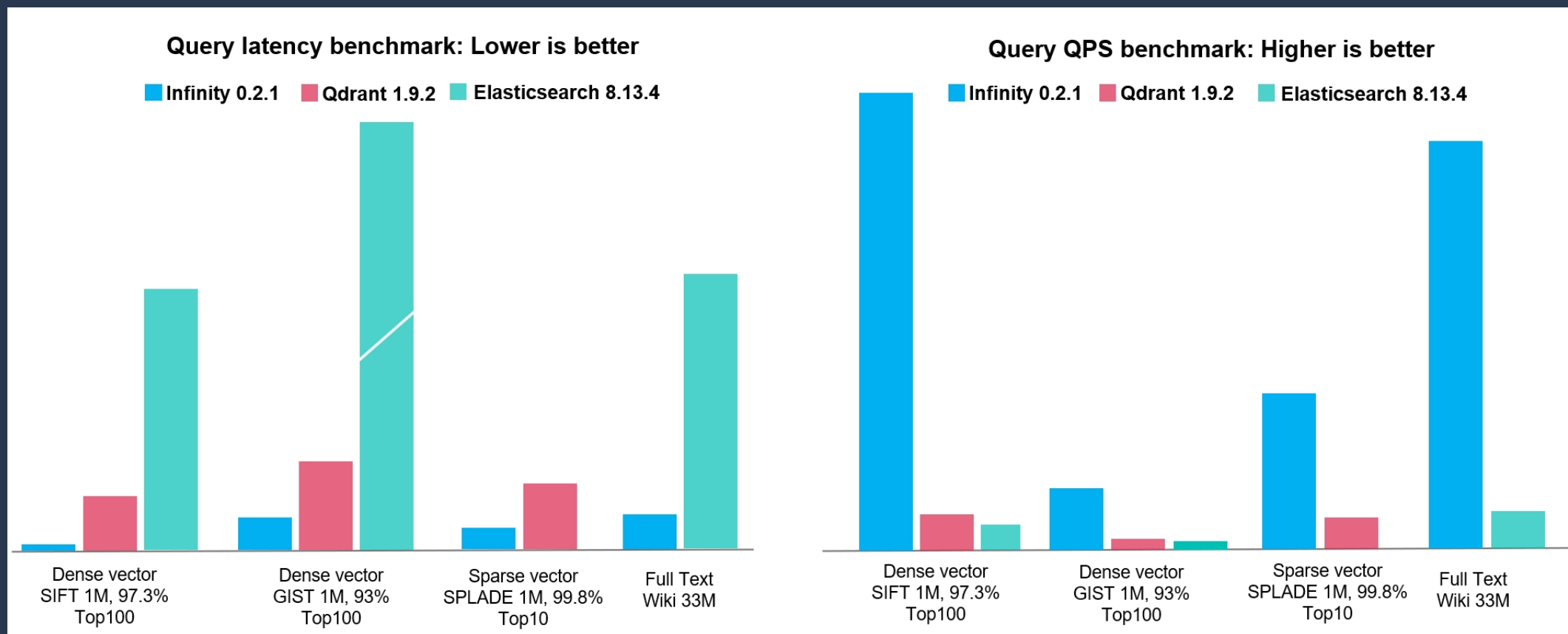
# 03

## 混合搜索

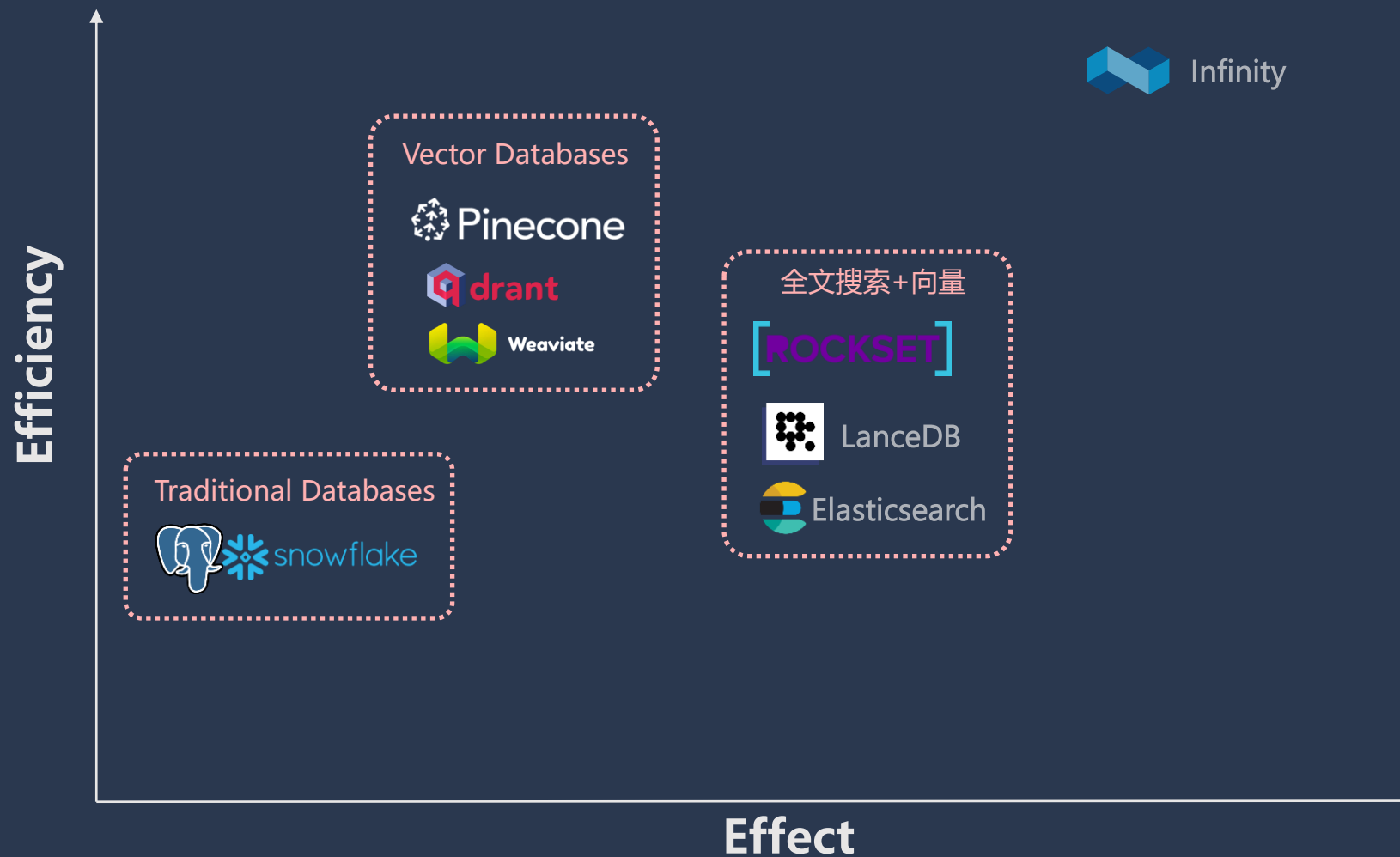
# Indexing Database



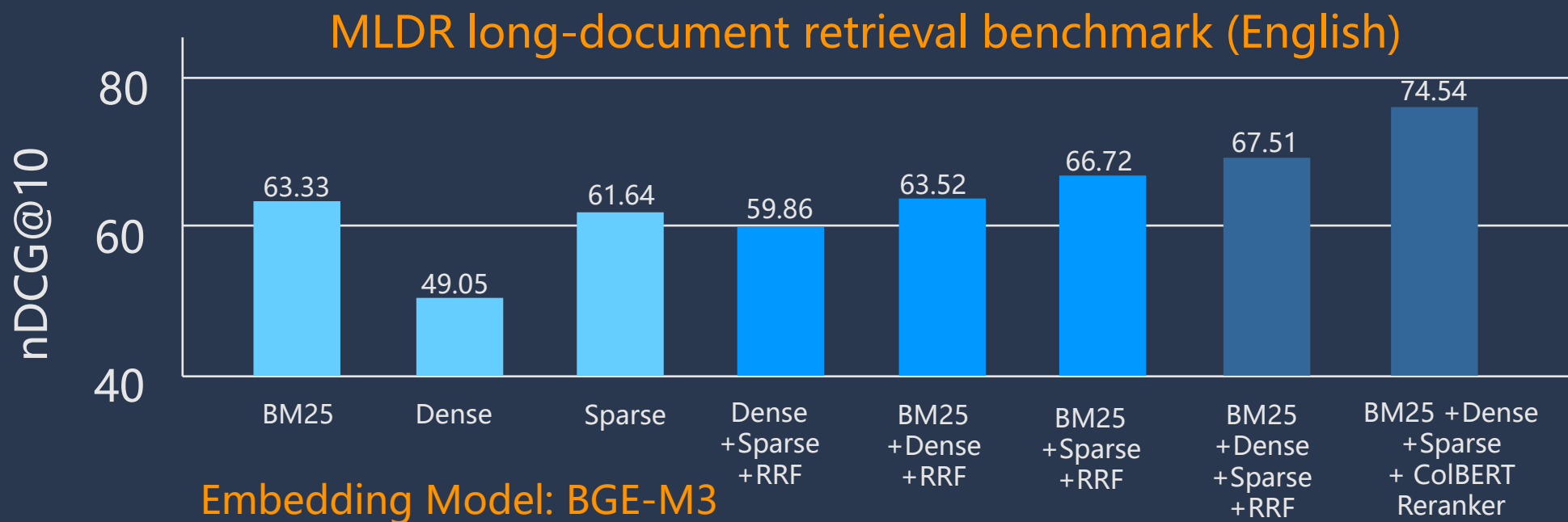
# Benchmark



# RAG数据库选型对比

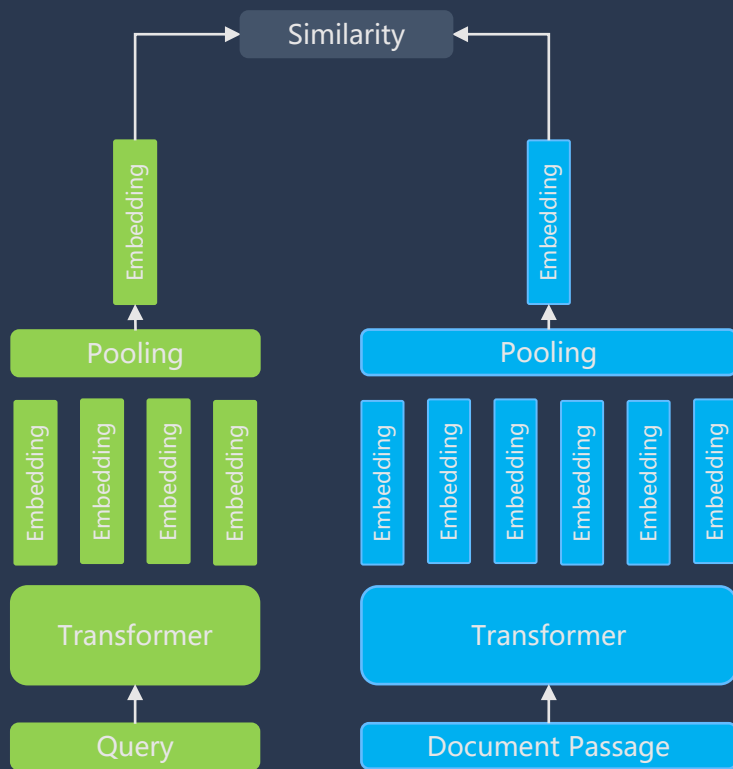


# 几路召回?

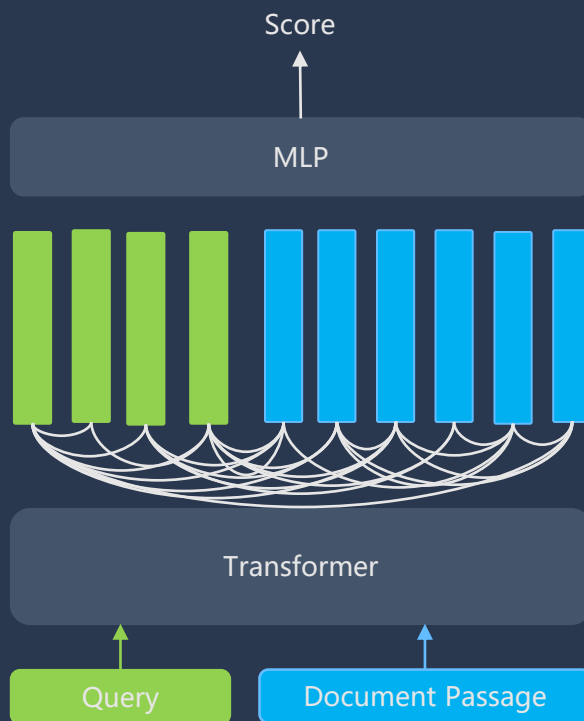


# 排序模型

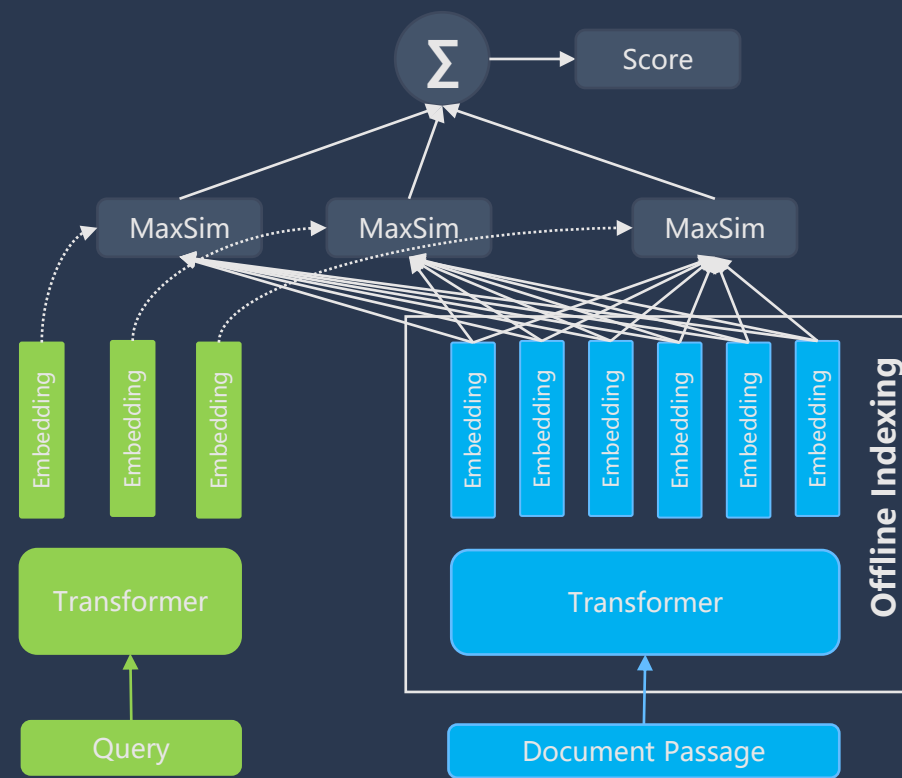
## Dual Encoder



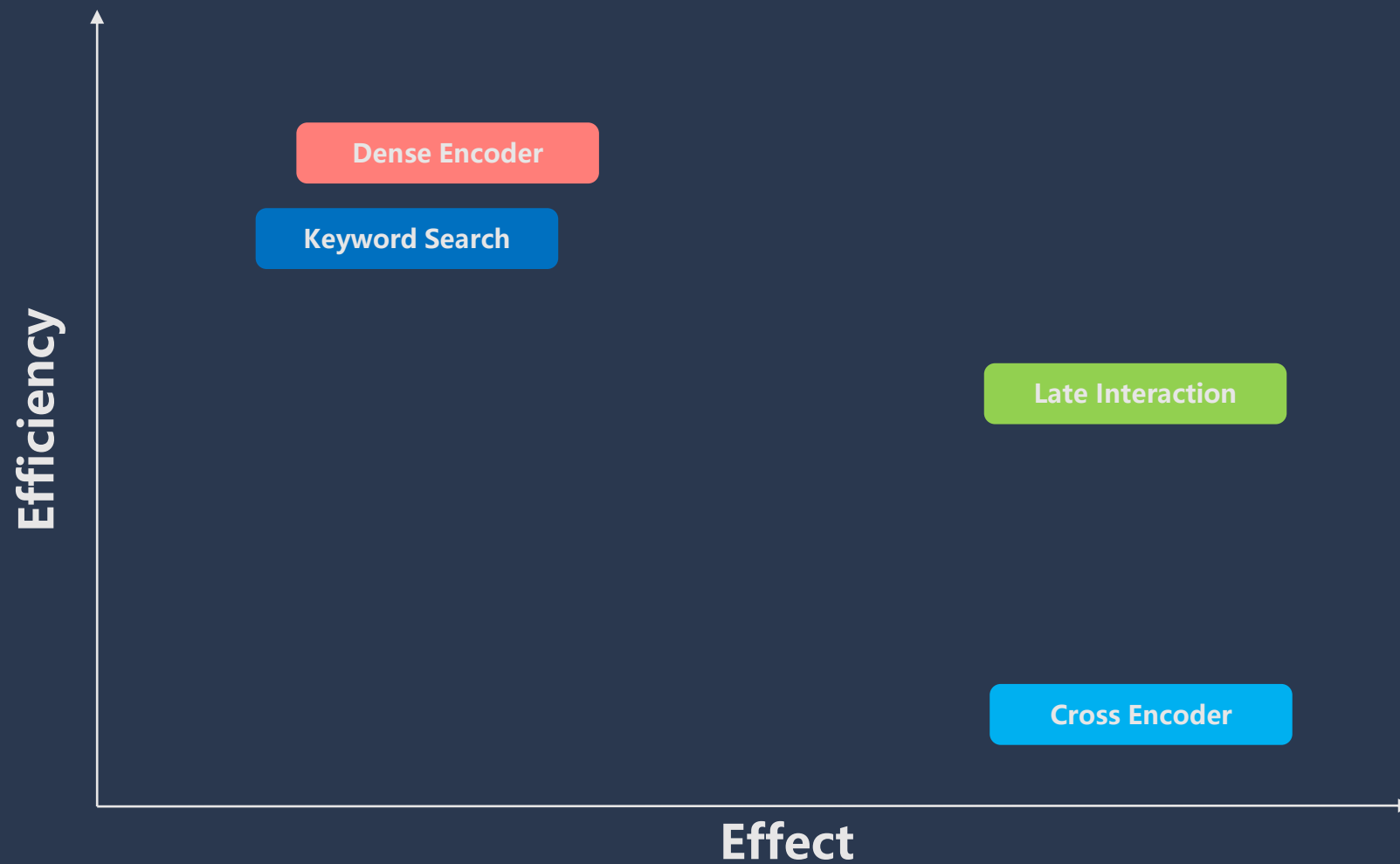
## Cross Encoder



## Late Interaction Encoder

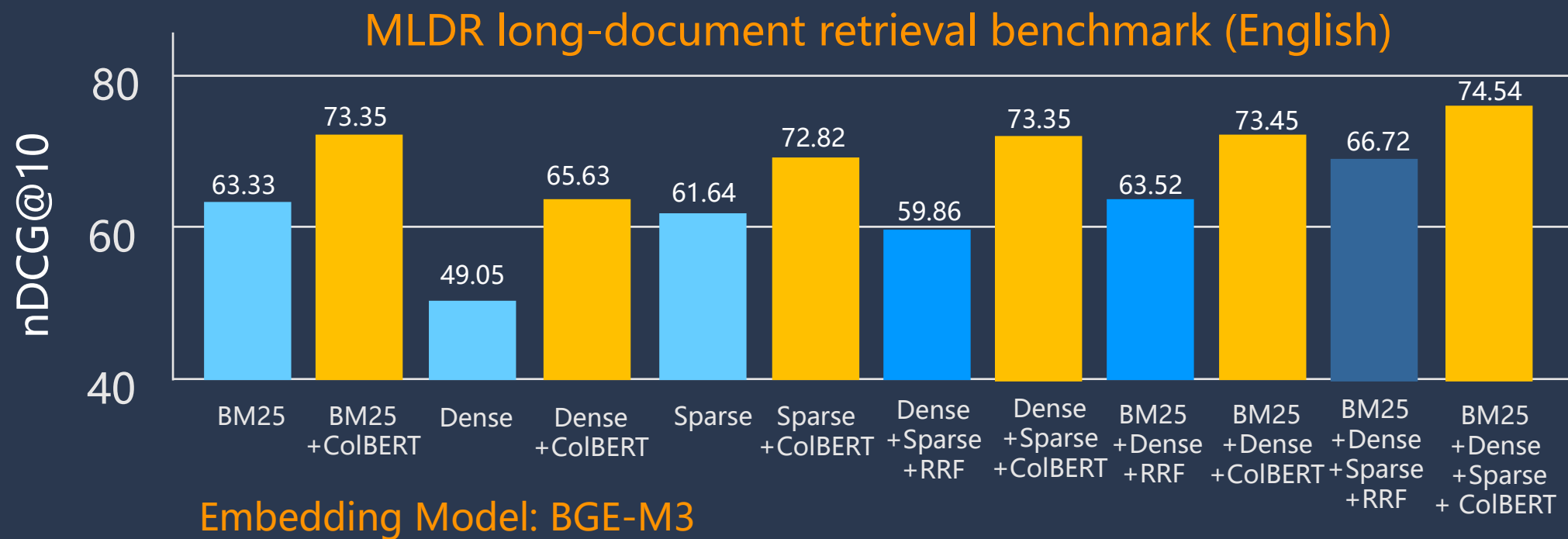


# ColBERT的收益

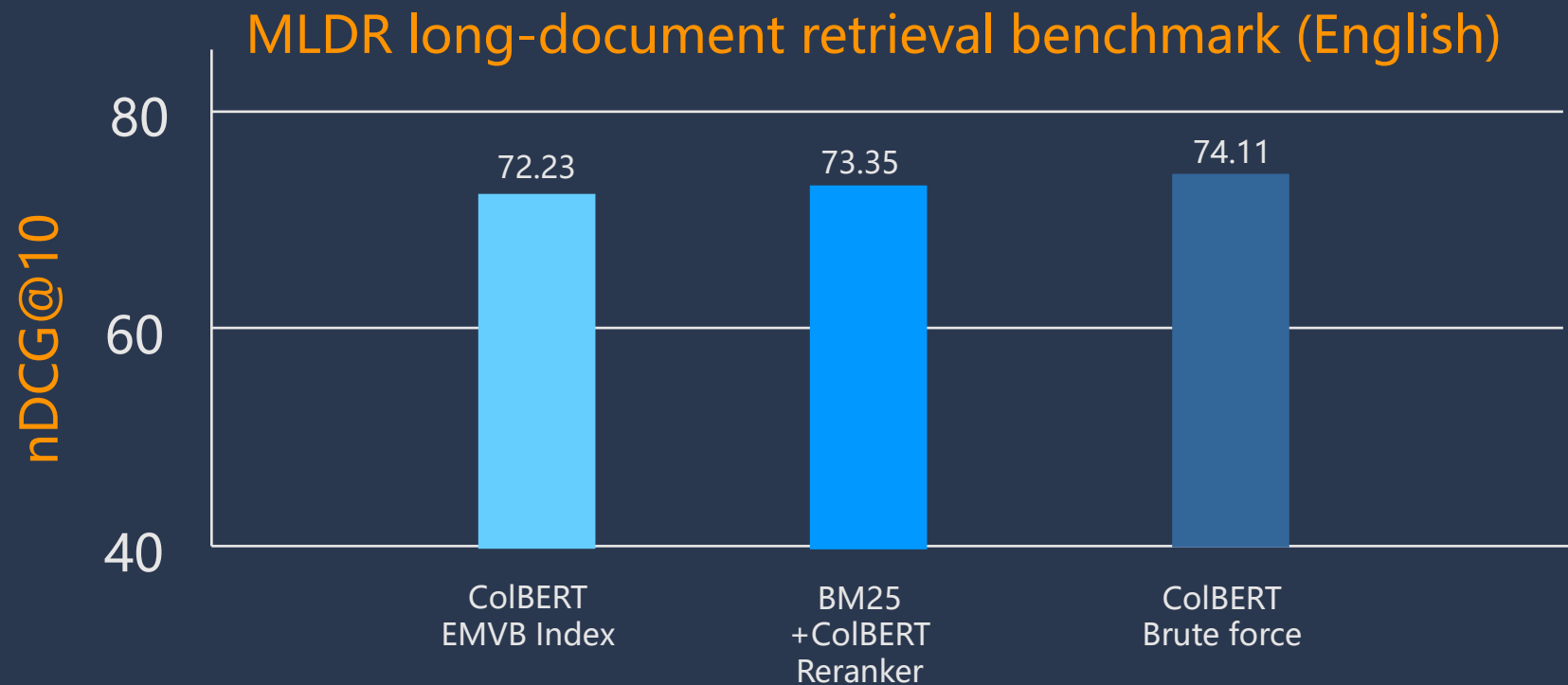




# ColBERT的收益

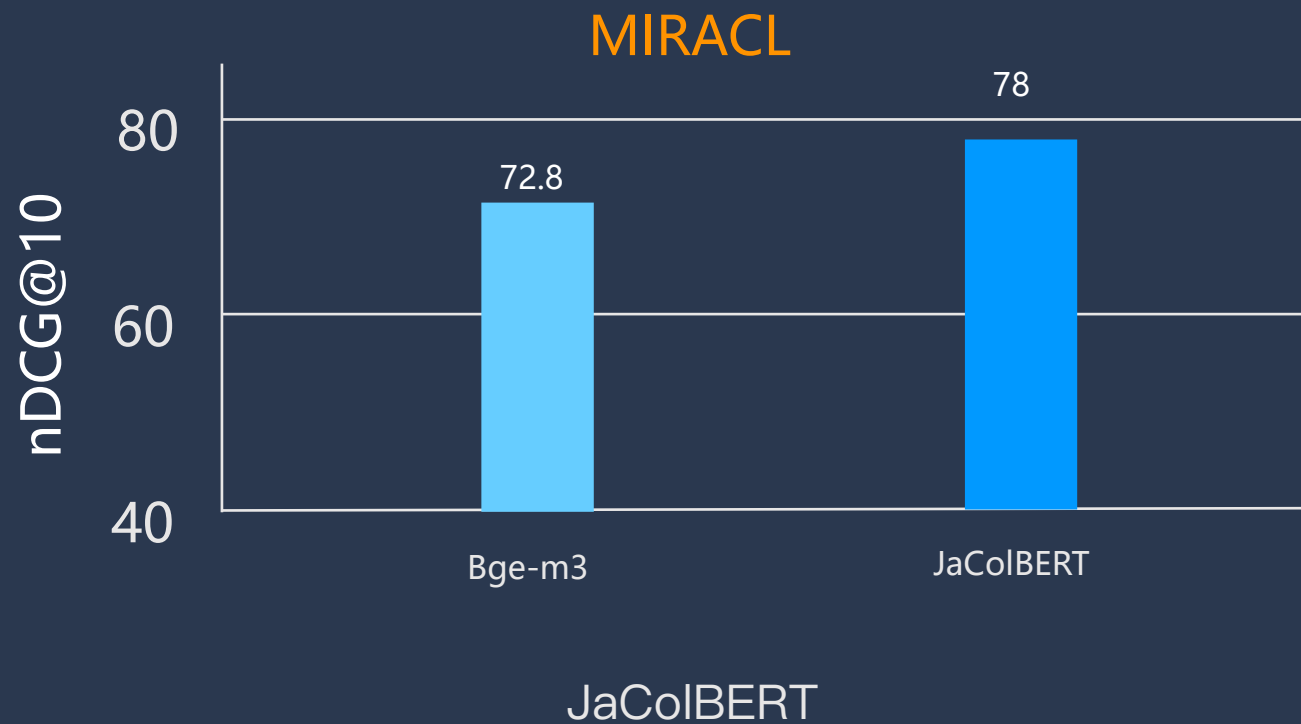


# ColBERT ranker 还是 reranker ?



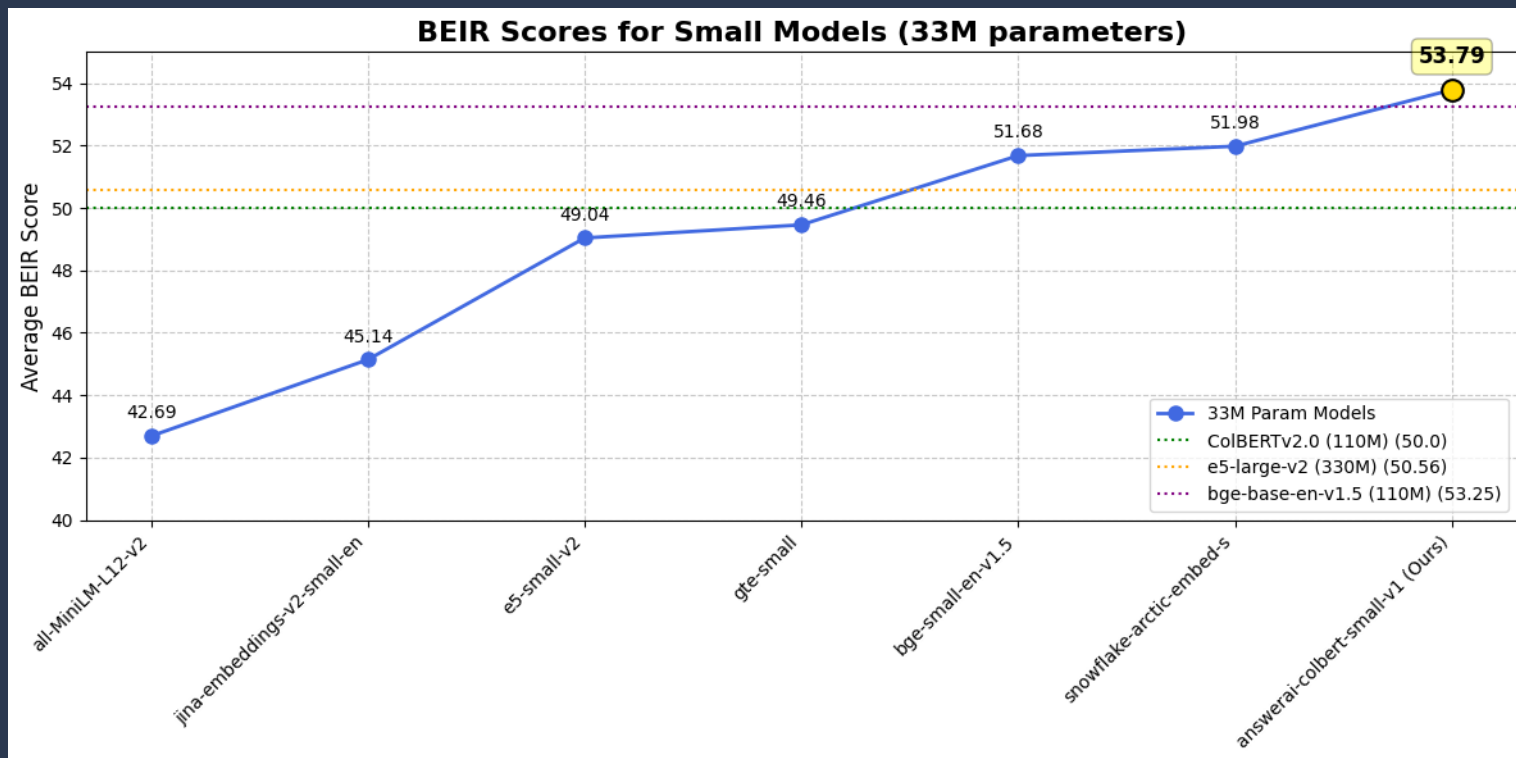
Embedding Model: BGE-M3

# 延迟交互是 RAG 的未来



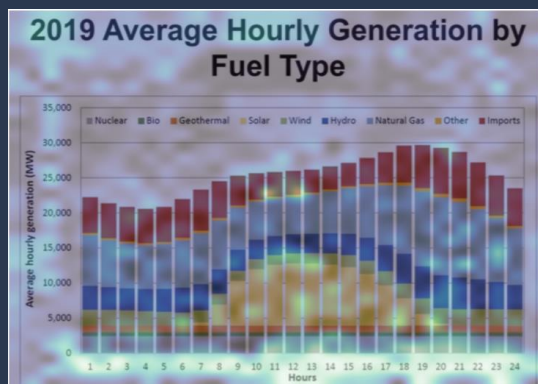
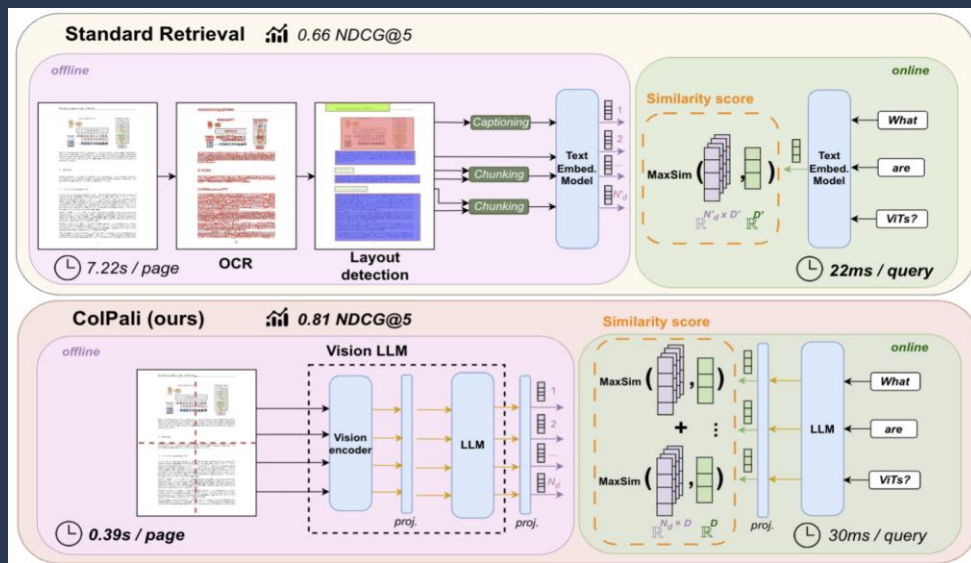
# 延迟交互是 RAG的未来

answerai-colbert-small-v1 基于JaCoBERT 33M参数



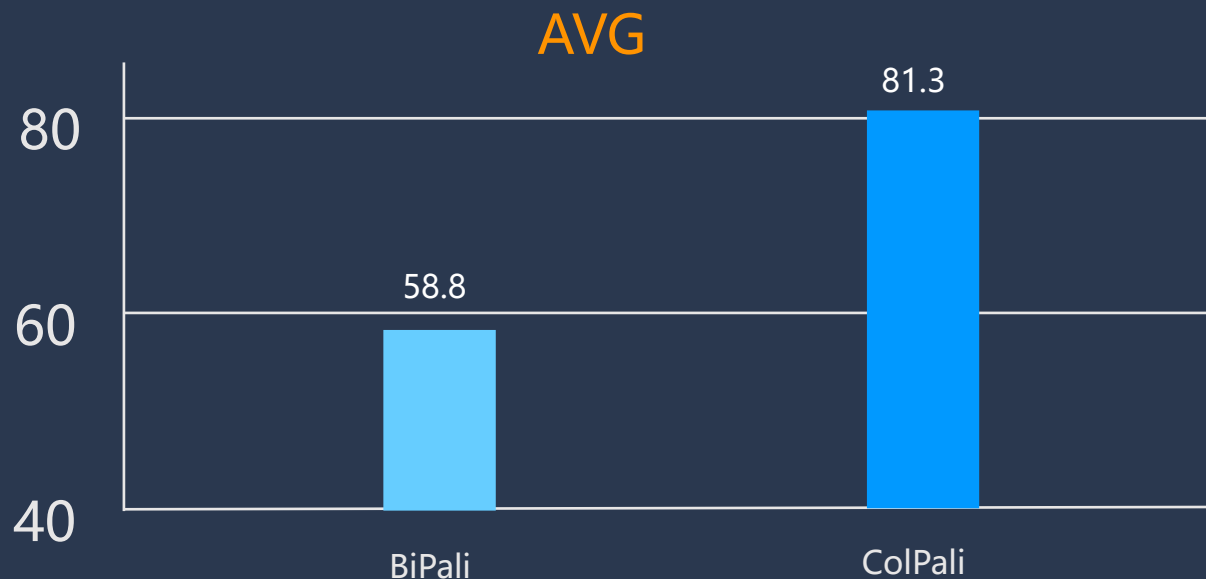
- 超过 BGE 110M
- 每个Token 96维
- Binary量化后每个Token 12 byte

# 延迟交互是 RAG 的未来



Query: Which hour of the day had the highest overall electricity generation in 2019 ?

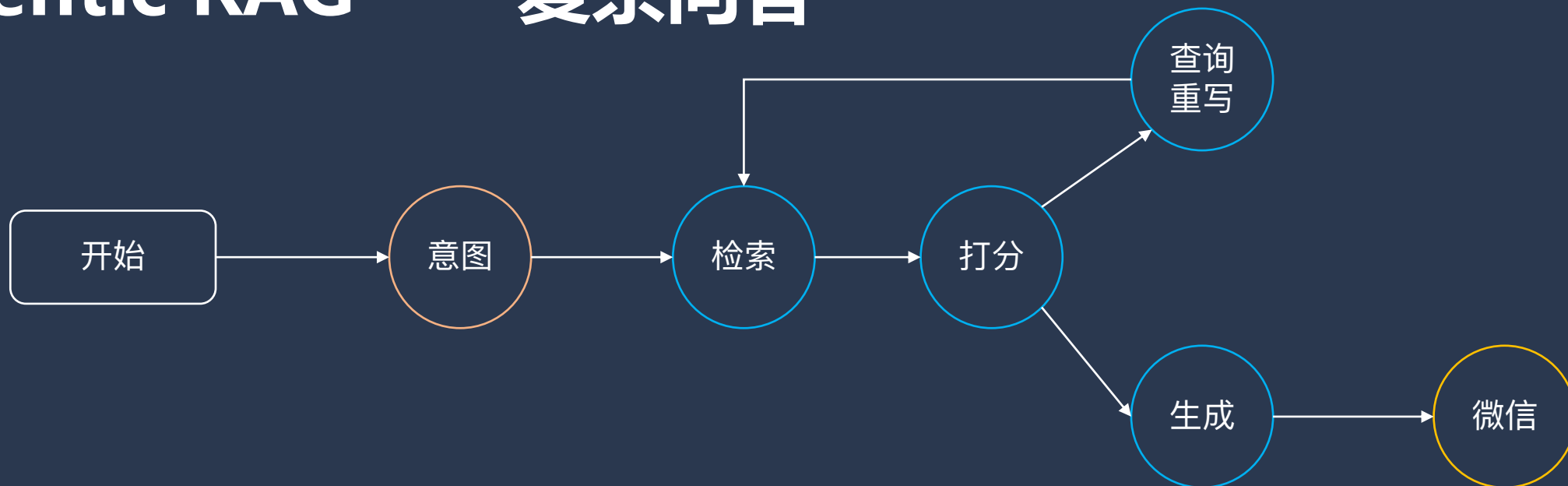
nDCG@5



ColPali

# 04 高级RAG

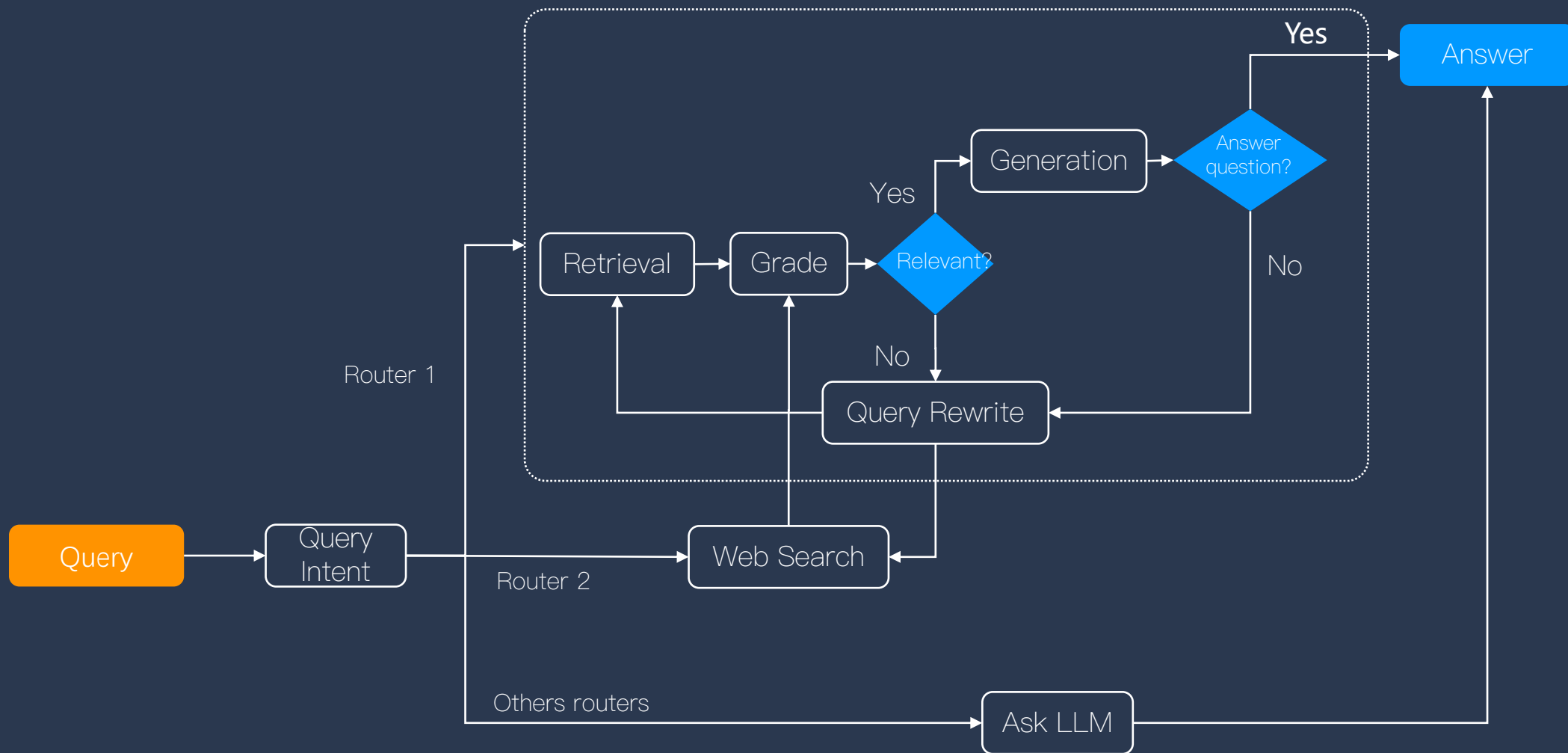
# Agentic RAG——复杂问答



Agentic RAG	
<ul style="list-style-type: none"><li>✓ 反思——自我纠错和迭代</li><li>✓ 工具——workflow</li><li>✓ 规划——workflow</li><li>✗ 多Agent协作</li></ul>	<ul style="list-style-type: none"><li>■ 查询意图识别</li><li>■ 知识图谱</li><li>■ 查询改写</li><li>■ 文档解析算法</li><li>■ 文档聚类 and 摘要RAPTOR</li><li>■ ...</li></ul>



# Agentic RAG——复杂问答



# 知识图谱

Data

Entities

Graph Construction and Augmentation

Passage

**Triplex**

Entity

Entity

Passage

Entity

Query

PPR

Passage

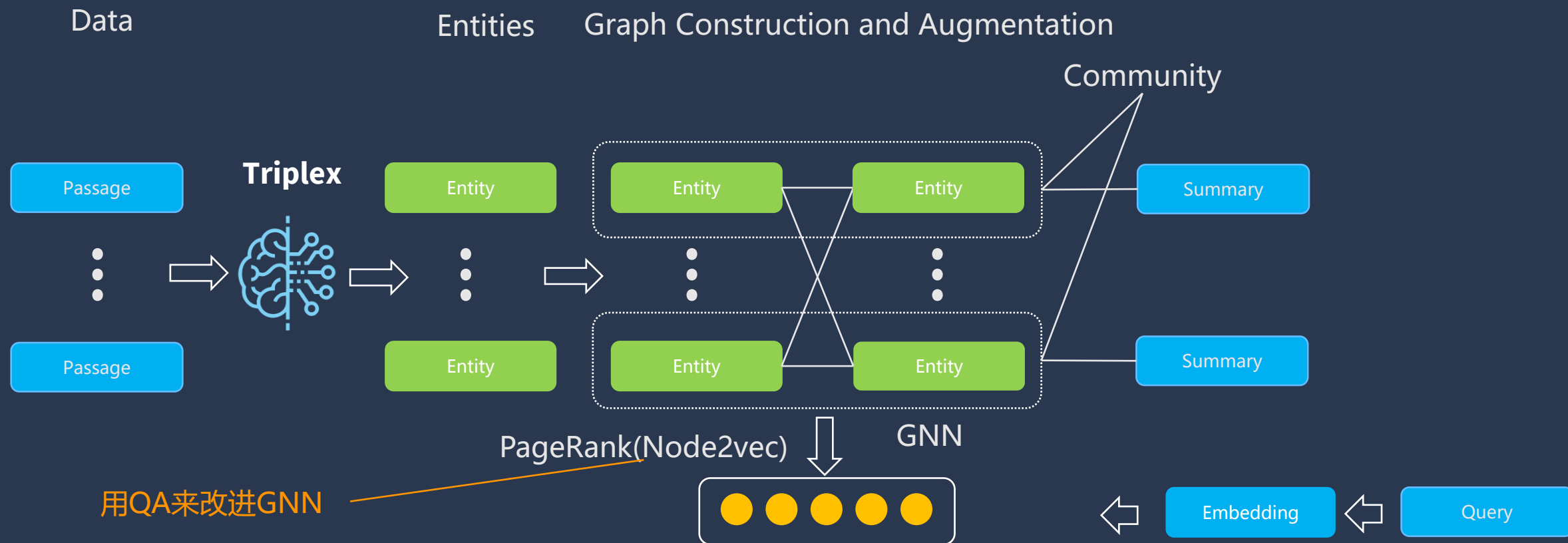
Entity

Entity

Passage

Entity

# 知识图谱



# 下一代RAG平台



# 极客邦科技 2024 年会议规划

促进软件开发及相关领域知识与创新的传播



# THANKS

智能未来，探索 AI 无限可能

Intelligent Future, Exploring the  
Boundless Possibilities of AI



<https://github.com/infiniflow/ragflow>  
<https://github.com/infiniflow/infinity>

**AiCon**  
全球人工智能开发与应用大会